



Computer  
Science

# **CSC196: Analyzing Data**

**Bayesian Probability and Inference**

Jason Pacheco and Cesim Erten

# Outline

- Introduction to Bayesian Probability
- Bayes Estimators
- Prediction and Updating
- Model Validation

# Outline

- Introduction to Bayesian Probability
- Bayes Estimators
- Prediction and Updating
- Model Validation

# What is Probability?

*What does it mean that the probability of heads is  $\frac{1}{2}$  ?*



*Two schools of thought...*

## **Frequentist Perspective**

Proportion of successes (heads) in repeated trials (coin tosses)

## **Bayesian Perspective**

Belief of outcomes based on assumptions about nature and the physics of coin flips

*Neither is better/worse, but we can compare interpretations...*

# Frequentist & Bayesian Modeling

*We will use the following notation throughout:*

$\theta$  - Unknown (e.g. coin bias)

$y$  - Data

## Frequentist

$$p(y; \theta)$$

- $\theta$  is a non-random unknown parameter
- $p(y; \theta)$  is the *sampling / data generating distribution*

## Bayesian

Prior Belief  $\rightarrow p(\theta)p(y | \theta) \leftarrow$  Likelihood

- $\theta$  is a random variable (latent)
- Requires specifying  $p(\theta)$  the prior belief

# Bayes' Rule

*Posterior represents all uncertainty after observing data...*

The diagram illustrates Bayes' Rule with the following components and labels:

- prior probability**: Labeled in red, with an arrow pointing to the  $p(\theta)$  term in the numerator.
- likelihood function for the parameters**: Labeled in red, with an arrow pointing to the  $p(y | \theta)$  term in the numerator.
- posterior probability**: Labeled in red, with an arrow pointing to the  $p(\theta | y)$  term on the left side of the equation.
- Marginal likelihood or: normalizer**: Labeled in red, with an arrow pointing to the  $p(y)$  term in the denominator.

$$p(\theta | y) = \frac{p(\theta)p(y | \theta)}{p(y)}$$

# Bayes' Rule : The Book's Notation

*Posterior represents all uncertainty after observing data...*

**prior** probability

**likelihood** function for the parameters

$$\pi(\theta | x) = \frac{\pi(\theta) f(x | \theta)}{g(x)}$$

**posterior** probability

**Marginal likelihood** or: **normalizer**

The diagram shows the equation for Bayes' Rule. Arrows point from the labels to the corresponding terms: 'prior probability' to  $\pi(\theta)$ , 'likelihood function for the parameters' to  $f(x | \theta)$ , 'posterior probability' to  $\pi(\theta | x)$ , and 'Marginal likelihood or: normalizer' to  $g(x)$ .

# Priors in AI / ML / Data Science

- Priors are often used as *regularizers* (promote smoothing)
  - Reduces overfitting as random noise is not smooth
  - Often regularizers can be of simple form, even conjugate
- Priors often house sophisticated domain knowledge
  - Possibly from earlier encounters with data
  - Possibly problem constraints (e.g.  $\theta$  must be nonnegative)
  - World knowledge is complex, so good priors are often complex and **not conjugate**

# Bayesian Inference Example

About **29%** of American adults have high blood pressure (BP). Home test has **30% false positive** rate and **no false negative error**.



A recent home test states that you have high BP. Should you start medication?

An Assessment of the Accuracy of Home Blood Pressure Monitors When Used in Device Owners

Jennifer S. Ringrose,<sup>1</sup> Gina Polley,<sup>1</sup> Donna McLean,<sup>2-4</sup> Ann Thompson,<sup>1,5</sup> Fraulein Morales,<sup>1</sup> and Raj Padwal<sup>1,4,6</sup>

# Bayesian Inference Example

About **29%** of American adults have high blood pressure (BP). Home test has **30% false positive** rate and **no false negative error**.



- Latent quantity of interest is hypertension:  $\theta \in \{true, false\}$
- Measurement of hypertension:  $y \in \{true, false\}$
- Prior:  $p(\theta = true) = 0.29$
- Likelihood:  $p(y = true \mid \theta = false) = 0.30$   
 $p(y = true \mid \theta = true) = 1.00$

# Bayesian Inference Example

About **29%** of American adults have high blood pressure (BP). Home test has **30% false positive** rate and **no false negative error**.



Suppose we get a positive measurement, then posterior is:

$$\begin{aligned} p(\theta = \text{true} \mid y = \text{true}) &= \frac{p(\theta = \text{true})p(y = \text{true} \mid \theta = \text{true})}{p(y = \text{true})} \\ &= \frac{0.29 * 1.00}{0.29 * 1.00 + 0.71 * 0.30} \approx 0.58 \end{aligned}$$

**What conclusions can be drawn from this calculation?**

# Bayes' Rule

*Posterior represents all uncertainty after observing data...*

The diagram shows the Bayes' Rule equation with four labels and arrows pointing to the corresponding parts of the equation:

- prior** probability: points to  $p(\theta)$
- likelihood** function for the parameters: points to  $p(y | \theta)$
- posterior** probability: points to  $p(\theta | y)$
- Marginal likelihood** or: **normalizer**: points to  $p(y)$

$$p(\theta | y) = \frac{p(\theta)p(y | \theta)}{p(y)}$$

# Bayes' Rule : Marginal Likelihood

$$p(\theta | y) = \frac{p(\theta)p(y | \theta)}{p(y)} \propto \underbrace{p(\theta)p(y | \theta)}$$

**Often hard to calculate**

**Often know this (the model)**

Marginal likelihood integrates (marginalizes) over unknown  $\theta$  :

$$p(y) = \int p(\theta)p(y | \theta) d\theta$$

**Marginal likelihood is less problematic in discrete models (not always)**

This integral often lacks a closed form and cannot be computed...

# Aside : Proportionality

Recall PMF / PDF must sum / integrate to 1,

$$\begin{array}{cc} \text{PMF} & \text{PDF} \\ \sum_x p(x) = 1 & \int p(x) dx = 1 \end{array}$$

May only know distribution constant that does not depend on RV  $x$ ,

$$\int \tilde{p}(x) dx = \mathcal{Z} \quad \text{so} \quad p(x) \propto \tilde{p}(x)$$

Properly normalized distribution by dividing our normalization constant:

$$\int p(x) dx = \int \frac{1}{\mathcal{Z}} \tilde{p}(x) dx = \frac{1}{\int \tilde{p}(x) dx} \int \tilde{p}(x) dx = 1$$

# Aside : Proportionality

**Example** Let  $X$  be a Bernoulli RV (coinflip) with probabilities *proportional to*:

$$\tilde{p}(X = 0) = 0.5$$

$$\tilde{p}(X = 1) = 1.5$$

Greater than 1, but  
It is an *unnormalized*  
probability

Compute normalization constant,

$$\mathcal{Z} = \tilde{p}(X = 0) + \tilde{p}(X = 1) = 2.0$$

Normalize probability distribution,

$$p(X) = \frac{1}{\mathcal{Z}} \tilde{p}(X) = \begin{pmatrix} 1/4 \\ 3/4 \end{pmatrix}$$

Sums to 1

# Outline

- Introduction to Bayesian Probability
- **Bayes Estimators**
- Prediction and Updating
- Model Validation

# Minimum Mean Squared Error (MMSE)

Posterior mean minimizes squared error,

$$\hat{\theta}^{\text{MMSE}} = \arg \min \mathbb{E}[(\hat{\theta} - \theta)^2 \mid y] = E[\theta \mid y]$$

- Minimizes error conditioned on observed data
- MMSE is an **unbiased estimator**
- MMSE is **asymptotically unbiased** and **asymptotically normal**,

$$\sqrt{N}(\hat{\theta}^{\text{MMSE}} - \theta) \rightarrow \mathcal{N}(0, \sigma^2)$$

# Example: Beta-Bernoulli MMSE

Let  $X_1, \dots, X_N \sim \text{Bernoulli}(\theta)$  and  $\theta \sim \text{Beta}(\alpha, \beta)$ .

- Beta is a distribution on probabilities  $\theta \in [0, 1]$
- *Shape* parameters  $\alpha$  and  $\beta$  with mean,

$$\mathbf{E}[\theta] = \frac{\alpha}{\alpha + \beta}$$

- Beta-Bernoulli has Beta posterior distribution,

$$p(\theta \mid X_1^N) = \text{Beta}(\alpha + \text{number of heads}, \beta + \text{number of tails})$$

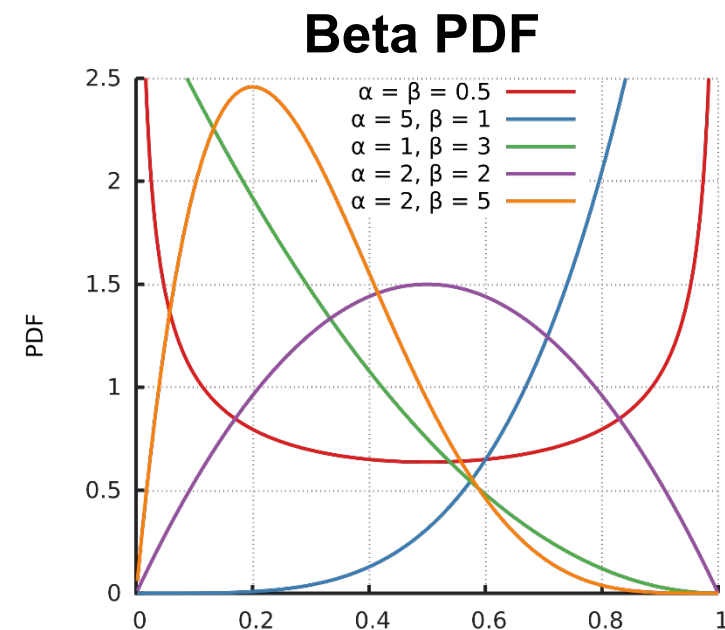
MMSE given by posterior mean,

**Q** What happens to MMSE when we have limited data?

$$\hat{\theta}^{\text{MMSE}} = \frac{\alpha + \text{number of heads}}{\alpha + \beta + N}$$

**Prior belief (pseudo-heads)**

**Q** What happens to MMSE when we have a lot of data?



# Bayes Estimators

Minimizes expected loss function,

$$\hat{\theta} = \arg \min_{\hat{\theta}} \mathbf{E} \left[ L(\theta, \hat{\theta}) \mid y \right]$$

Expected loss referred to as *Bayes risk*.

**MMSE** minimizes squared-error loss  $L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$

**Minimum absolute error (MAE)** is posterior *median*,

$$\arg \min \mathbf{E}[|\hat{\theta} - \theta| \mid y] = \text{median}(\theta \mid y)$$

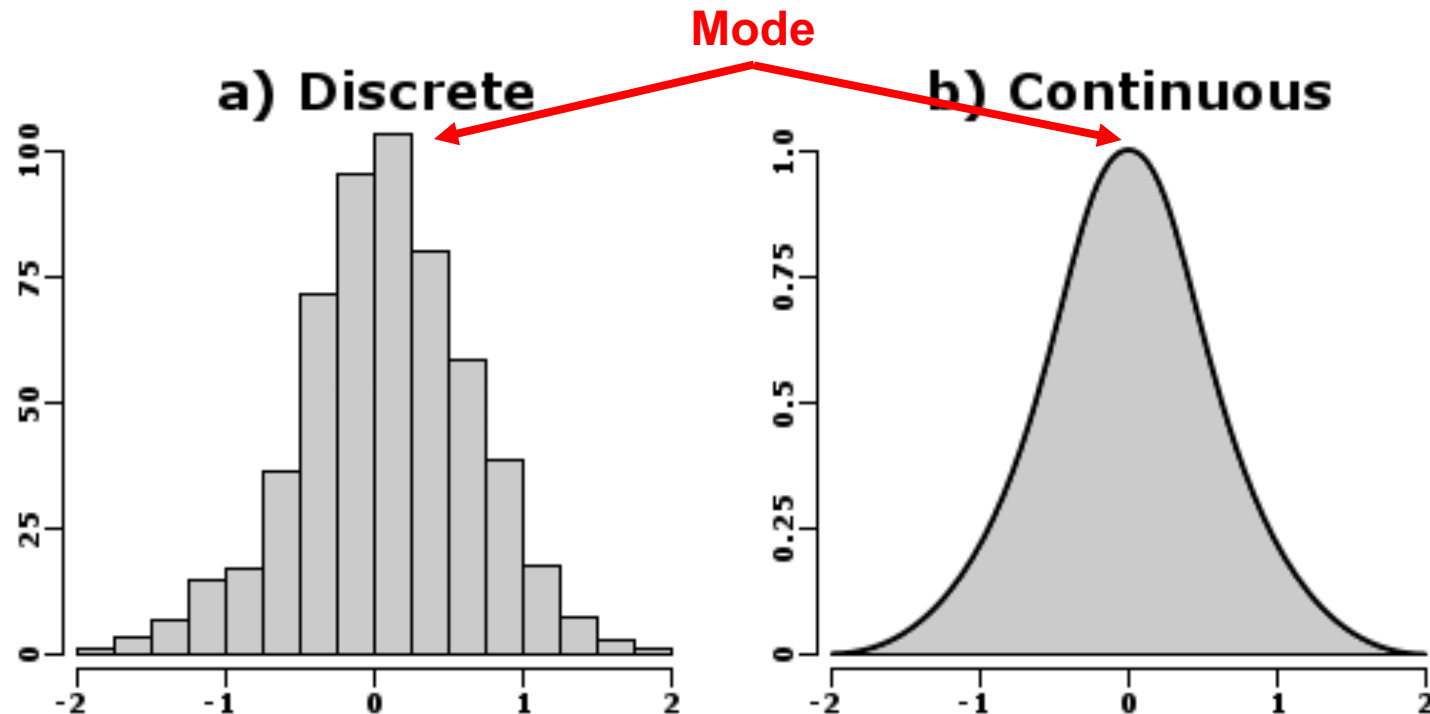
Note: Same answer for linear function:  $L(\theta, \hat{\theta}) = c|\hat{\theta} - \theta|$

# Maximum a Posteriori (MAP)

Very common to produce maximum probability estimates,

$$\hat{\theta}^{\text{MAP}} = \arg \max p(\theta | y)$$

*MAP is the **mode** ( highest probability outcome ) of the posterior*



# Example: Beta-Bernoulli MAP

Let  $X_1, \dots, X_N \sim \text{Bernoulli}(\theta)$  and  $\theta \sim \text{Beta}(\alpha, \beta)$  then posterior is,

$$p(\theta | X_1^N) = \text{Beta}(\alpha + \underbrace{\text{number of heads}}_{N_H}, \beta + \text{number of tails})$$

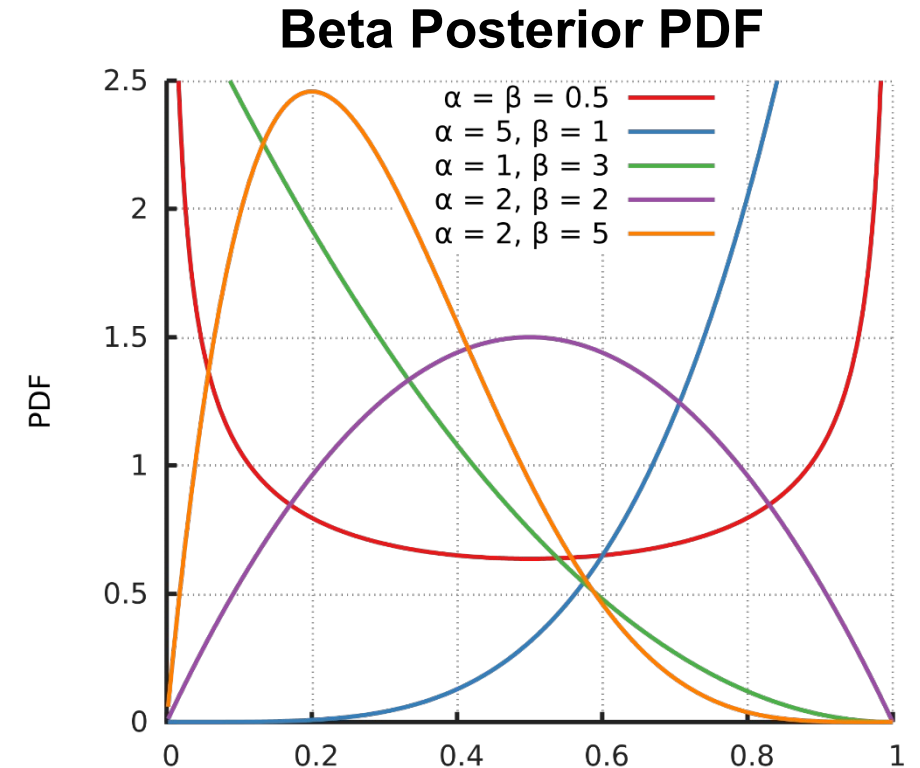
Highest probability (mode) of Beta given by,

Take derivative,  
set to zero, solve.

$$\hat{\theta}^{\text{MAP}} = \frac{\alpha + N_H - 1}{\alpha + \beta + N - 2}$$

Beta distribution is not always convex!

- MAP is any value for  $\alpha = \beta = 1$
- Two modes (bimodal) for  $\alpha, \beta < 1$



# Maximum a Posteriori (MAP)

Equivalent to maximizing joint probability,

$$\arg \max_{\theta} p(\theta | y) = \arg \max_{\theta} \frac{p(\theta, y)}{p(y)} = \arg \max_{\theta} p(\theta, y)$$

**Constant**

For iid  $y_1, \dots, y_N$  solve in log-domain (like *maximum likelihood est.*),

$$\hat{\theta}^{\text{MAP}} = \arg \max_{\theta} \log p(\theta, y_1, \dots, y_N) = \underbrace{\sum_i \log p(y_i | \theta)}_{\text{Log-Likelihood (how well it fits data)}} + \underbrace{\log p(\theta)}_{\text{Log-Prior (how well it agrees with prior)}}$$

***Intuition*** MAP is like MLE but with a “penalty” term (log-prior)

# Outline

- Introduction to Bayesian Probability
- Bayes Estimators
- **Prediction and Updating**
- Model Validation

# Prediction

Can make predictions of unobserved  $\tilde{y}$  before seeing any data,

$$p(\tilde{y}) = \sum_k p(\theta = k)p(\tilde{y} | \theta = k)$$

**Similar calculation to marginal likelihood**

*This is the **prior predictive** distribution*

For continuous parameters sum turns into integral,

$$p(\tilde{y}) = \int p(\theta)p(\tilde{y} | \theta) d\theta$$

*This is a prediction based on **no observed data***

# Prediction

When we observe  $y$  we can predict future observations  $\tilde{y}$ ,

$$p(\tilde{y} | y) = \sum_k \underbrace{p(\theta = k | y)}_{\text{This is now the posterior}} p(\tilde{y} | \theta = k)$$

This is now the posterior

*This is the **posterior predictive distribution***

Again, for continuous parameters sum turns into integral,

$$p(\tilde{y} | y) = \int p(\theta | y) p(\tilde{y} | \theta) d\theta$$

# Prediction Example

About **29%** of American adults have high blood pressure (BP). Home test has **30% false positive** rate and no false negative error.



What is the likelihood of *another* positive measurement?

$$p(\tilde{y} = true \mid y = true) = \sum_{\theta \in \{true, false\}} p(\theta \mid y = true) p(\tilde{y} = true \mid \theta)$$

$$= 0.42 * 0.30 + 0.58 * 1.00 \approx 0.71$$

**What conclusions can be drawn from this calculation?**

# Bayesian Updating

*Suppose we plan to take another test...*

**Question** What is our belief about blood pressure status *before* the second test?

(a) Posterior:  $p(\theta = \text{true} \mid y_1 = \text{true})$

(b) Likelihood:  $p(y_1 = \text{true} \mid \theta = \text{true})$

(c) Marginal Likelihood:  $p(y_1 = \text{true})$

# Bayesian Updating

*Suppose we plan to take another test...*

**Question** What is the probability that we get *true* on the second test if we have high blood pressure?

(a) Posterior:  $p(\theta = \text{true} \mid y_1 = \text{true}, y_2 = \text{true})$

(b) Likelihood:  $p(y_2 = \text{true} \mid \theta = \text{true})$

(c) Marginal Likelihood:  $p(y_2 = \text{true})$

*Why not:*  $p(y_2 = \text{true} \mid \theta = \text{true}, y_1 = \text{true})$

# Bayesian Updating

*Suppose we plan to take another test...*

**Question** What is the probability that we get *true* on the second test if we have high blood pressure?

(a) Posterior:  $p(\theta = \text{true} \mid y_1 = \text{true}, y_2 = \text{true})$

(b) Likelihood:  $p(y_2 = \text{true} \mid \theta = \text{true})$

(c) Marginal Likelihood:  $p(y_2 = \text{true})$

*Because*  $y_1 \perp y_2 \mid \theta$  *so*  $p(y_2 \mid \theta, y_1) = p(y_2 \mid \theta)$

# Outline

- Introduction to Bayesian Probability
- Bayes Estimators
- Prediction and Updating
- **Model Validation**

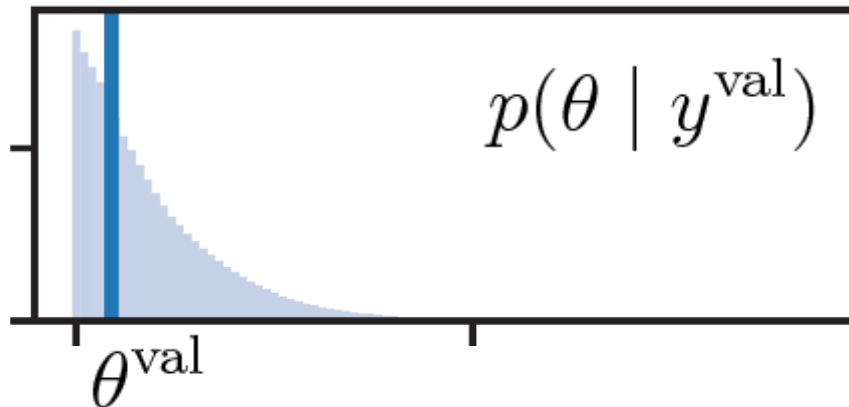
# Model Validation

*How do we know if the model  $p(\theta, y)$  is good?*

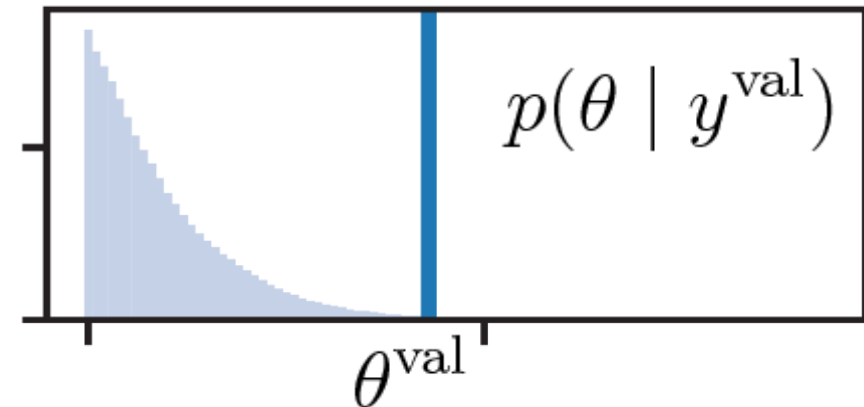
## Supervised Learning

Validation set  $\{(\theta^{\text{val}}, y^{\text{val}})\}$  consists of known  $\theta^{\text{val}}$ . Are true values typically preferred under the posterior?

Good (maybe lucky)



Not Good (maybe unlucky)



Repeat trials over validation set for more certainty

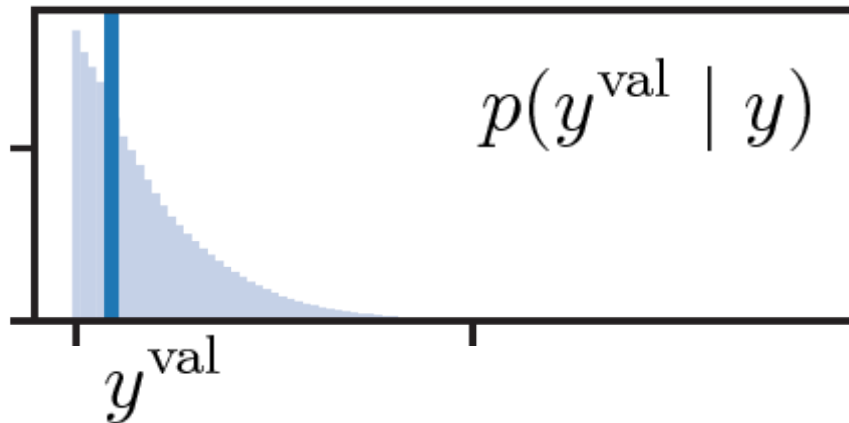
# Model Validation

*How do we know if the model  $p(\theta, y)$  is good?*

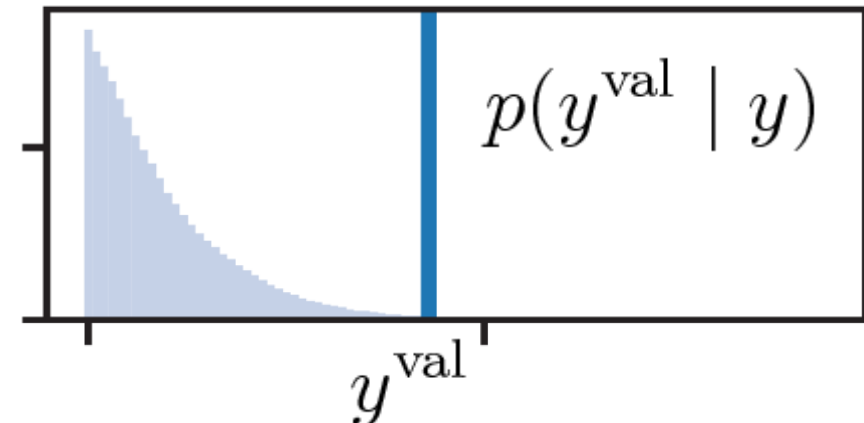
## Unsupervised Learning

Validation set  $\{y^{\text{val}}\}$  only contains observable data. Check validation data against posterior-predictive distribution.

Good (maybe lucky)



Not Good (maybe unlucky)



Repeat trials over validation set for more certainty

# Likelihood and Odds Ratios

Which parameter value  $\theta_1$  or  $\theta_2$  is more likely to have generated the observed data  $y$ ?

The **posterior odds ratio** is:

$$\frac{p(\theta_1 | y)}{p(\theta_2 | y)} = \frac{p(\theta_1) p(y | \theta_1) \cancel{p(y)}}{p(\theta_2) p(y | \theta_2) \cancel{p(y)}}$$

Prior Odds  
Ratio

Likelihood  
Ratio

**Observe:** the marginal likelihood  $p(y)$  cancels!

# Posterior Summarization

*Ideally we would report the full posterior distribution as the result of inference...but this is not always possible*

## **Summary of Posterior Location:**

Point estimates: mean (MMSE), mode, median (min. absolute error)

## **Summary of Posterior Uncertainty:**

Credible intervals / regions, posterior entropy, variance

**Bayesian analysis should report uncertainty when possible**

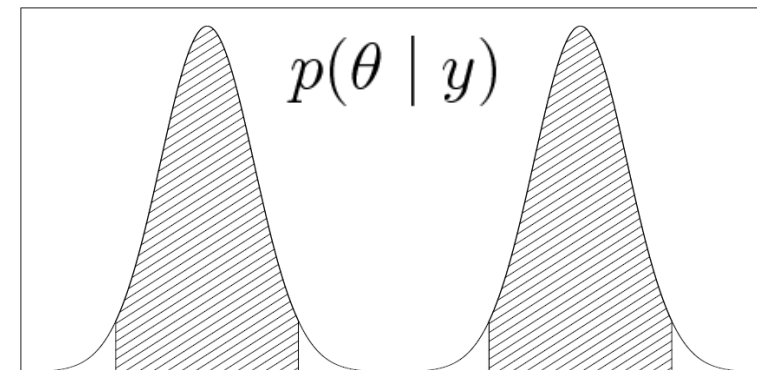
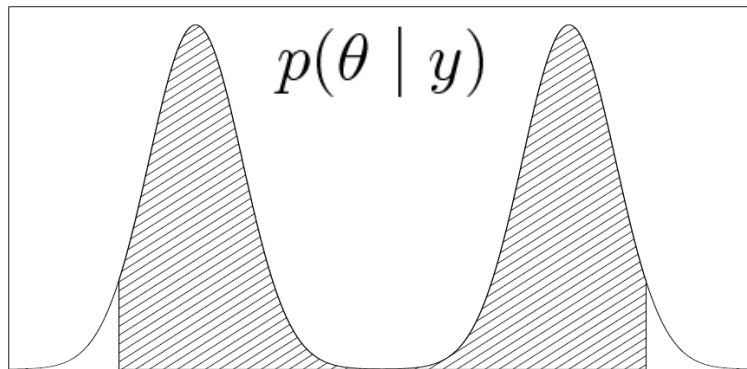
# Credible Interval

**Def.** For parameter  $0 < \alpha < 1$  the  $100(1 - \alpha)\%$  credible interval  $(L(y), U(y))$  satisfies,

$$p(L(y) < \theta < U(y) \mid y) = \int_{L(y)}^{U(y)} p(\theta \mid y) = 1 - \alpha$$

**Interval containing fixed percentage of posterior probability density.**

**Note:** This is not unique -- consider the 95% intervals below:



# Frequentist Inference

**Example:** Suppose we observe the outcome of  $N$  coin flips.  
 $y = \{y_1, \dots, y_N\}$ . What is the probability of heads  $\theta$  (coin bias)?

- Coin bias  $\theta$  is not random (e.g. there is some *true* value)
- Uncertainty reported as confidence interval (typically 95%)

Correct Interpretation: On repeated trials of  $N$  coin flips  $\theta$  will fall inside the confidence interval 95% of the time (in the limit)

- Inferences are valid for multiple trials, **never on single trials**

**Wrong Interpretation:** For *this trial* there is a 95% chance  $\theta$  falls in the confidence interval

# Bayesian Inference

*Posterior distribution is complete representation of uncertainty*

$$p(\theta | y) = \frac{p(\theta)p(y | \theta)}{p(y)}$$

**Prior Belief**      **Likelihood**  
**Marginal Likelihood**  
**(more on this later)**

- Must specify a prior belief  $p(\theta)$  about coin bias
- Coin bias  $\theta$  is a random quantity
- Interval  $p(l(y) < \theta < u(y) | y) = 0.95$  can be reported in lieu of full posterior, and takes intuitive interpretation for a single trial

Interval Interpretation: For this experiment there is a 95% chance that  $\theta$  lies in the interval

# Summary

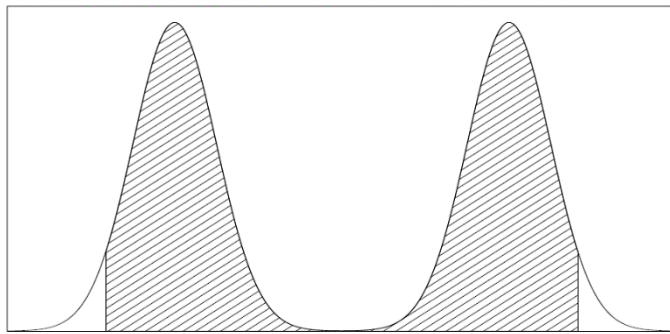
- Bayesian statistics interprets probability differently than classical stats
  - Frequentist: Probability  $\rightarrow$  Long run odds in repeated trials
  - Bayesian: Probability  $\rightarrow$  Belief of outcome that captures all uncertainty
- Bayesian models treat unknown parameter as random, with a prior
- Bayesian inference via the *posterior distribution* using Bayes' rule

$$p(\theta | y) = \frac{p(\theta)p(y | \theta)}{p(y)}$$

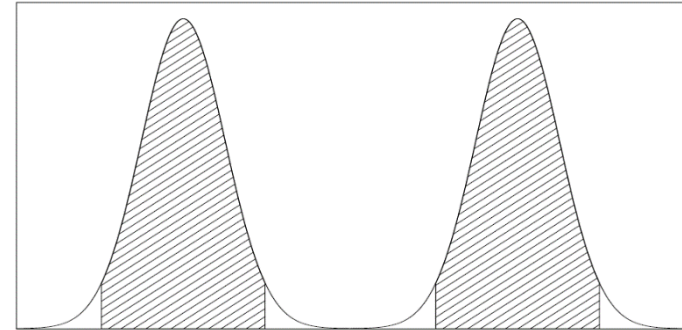
- Bayesian estimators minimize expected risk (e.g. MMSE)
- Maximum a posteriori (MAP) estimate maximizes posterior probability

# Summary

- Conjugate prior-posterior pairs ensure closed-form posterior inference
- Posterior uncertainty can be characterized by credible intervals



Not necessarily  
unique



- Selecting models can be done via posterior odds ratio

$$\frac{p(\theta_1 | y)}{p(\theta_2 | y)} = \frac{p(\theta_1) p(y | \theta_1) \cancel{p(y)}}{p(\theta_2) p(y | \theta_2) \cancel{p(y)}}$$

- Parameter can be marginalized out via prior/posterior predictive dist'n